

Robust and Label-Efficient Deep Waste Detection

Hassan Abid^{†,1}

hassan.abid@mbzuai.ac.ae

Khan Muhammad^{†,2}

khan.muhammad@ieee.org

Muhammad Haris Khan¹

muhammad.haris@mbzuai.ac.ae

¹ Mohamed Bin Zayed University of
Artificial Intelligence,
Abu Dhabi, UAE

² Sungkyunkwan University,
Seoul, South Korea

Abstract

Effective waste sorting is critical for sustainable recycling, yet AI research in this domain continues to lag behind commercial systems due to limited datasets and reliance on legacy object detectors. In this work, we advance AI-driven waste detection by establishing strong baselines and introducing an ensemble-based semi-supervised learning framework. We first benchmark state-of-the-art Open-Vocabulary Object Detection (OVOD) models on the real-world ZeroWaste dataset, demonstrating that while class-only prompts perform poorly, LLM-optimized prompts significantly enhance zero-shot accuracy. Next, to address domain-specific limitations, we fine-tune modern transformer-based detectors, achieving a new baseline of 51.6 mAP. We then propose a soft pseudo-labeling strategy that fuses ensemble predictions using spatial and consensus-aware weighting, enabling robust semi-supervised training. Applied to the unlabeled ZeroWaste-s subset, our pseudo-annotations achieve performance gains that surpass fully supervised training, underscoring the effectiveness of scalable annotation pipelines. Our work contributes to the research community by establishing rigorous baselines, introducing a robust ensemble-based pseudo-labeling pipeline, generating high-quality annotations for the unlabeled ZeroWaste-s subset, and systematically evaluating OVOD models under real-world waste sorting conditions. Our code is available at: [GitHub Repository](#).

Introduction

The global waste crisis, driven by urbanization and increased consumption, poses a critical environmental, health, and economic challenge. The generation of municipal solid waste is projected to grow from 2.01 billion tonnes in 2016 to 3.40 billion tonnes by 2050 [19]. Alarmingly, only 13.5% of waste is recycled, while over 33% is improperly managed, causing severe global challenges [19, 63]. Improving recycling efficiency is essential for the advancement of several Sustainable Development Goals (SDGs) of the United Nations [20, 68]. Developing intelligent and scalable waste sorting technologies has therefore become a key priority for global sustainability efforts.

[†] Corresponding authors.

© 2025. The copyright of this document resides with its authors.

It may be distributed unchanged freely in print or electronic forms.

Material Recovery Facilities (MRFs) play a central role by sorting waste into recyclable streams such as plastics, metals, paper, and cardboard [1, 2]. However, despite utilizing machinery alongside manual labor [3], recycling rates remain low, with less than 35% of recyclable waste recovered in the United States as of 2018 [4]. Conventional MRF operations suffer from inefficiencies and high material losses, discarding up to 20% of recyclables [5], and expose workers to hazardous conditions including sharp objects, toxic substances, and medical waste [3]. Moreover, the cluttered, overlapping, and deformable nature of waste streams makes automated detection and sorting particularly challenging. To address these limitations, recent advances have integrated Artificial Intelligence (AI) and robotics into MRF workflows, with object detection models emerging as a practical foundation for scalable, real-time sorting under complex conditions [6, 7]. Although several commercial companies [8, 9, 10] have developed effective AI-powered waste sorting solutions, their models and datasets remain proprietary, limiting broader research progress. In contrast, most academic studies rely on outdated object detection baselines and simplistic datasets collected in controlled environments [11, 12, 13, 14, 15], which fail to reflect the real-world challenges of industrial waste streams. Consequently, research in AI-driven waste detection significantly lags behind commercial advances, hindering reproducibility and scalability. Additionally, recent progress in OVOD [16, 17, 18, 19], which leverages vision-language models (VLMs) to generalize beyond fixed category sets, offers a promising avenue for improving adaptability. However, the effectiveness of OVOD models in complex industrial waste sorting scenarios remains largely unexplored.

In this work, we make four key contributions: (1) We develop an ensemble-based pseudo-labeling pipeline to generate high-quality annotations for the previously unlabeled ZeroWaste dataset, enabling scalable benchmarking and reducing reliance on costly manual labeling. (2) We conduct a comprehensive zero-shot evaluation of state-of-the-art OVOD models on real-world waste sorting data, highlighting their strengths and limitations in cluttered, deformable environments. (3) We design a robust semi-supervised learning framework that leverages ensemble-based pseudo-labels to substantially improve detection performance through large-scale training on unlabeled data. (4) We establish new strong baselines for waste detection by fine-tuning advanced object detectors, providing critical benchmarks to guide future research in AI-driven waste recovery.

2 Related Work

Deep Learning for Waste Recognition. Deep learning has been widely adopted for waste recognition tasks. Early works focused on image-level classification, utilizing convolutional neural networks (CNNs) trained on datasets such as TrashNet [20] and TACO [21]. RecycleNet [2] and Vo *et al.* [22] improved classification using DenseNet and ResNeXt architectures, while later studies explored lightweight models [23, 24], metadata integration [25], and attention mechanisms [26, 27]. However, these methods assume isolated objects in uncluttered settings, limiting their applicability to real-world MRFs characterized by severe clutter, deformation, and occlusion. To enable robotic sorting, research shifted to object detection frameworks capable of simultaneous classification and localization. Faster R-CNN [28, 29] and YOLO variants [30, 31, 32] were adapted for waste detection, with lightweight improvements for real-time inference [3, 33, 34, 35]. However, evaluations largely remained on controlled datasets. We address this limitation by benchmarking modern state-of-the-art detectors under realistic MRF conditions using the ZeroWaste dataset [2].

Waste Recognition Datasets. Public datasets for waste recognition vary widely in complexity and realism. Early datasets [21, 33, 32, 40, 50, 59] feature isolated objects, synthetic imagery, or simplified backgrounds that fail to capture the complexities of real-world settings. Recent datasets improve realism by collecting data in operational facilities [8, 60]. However, they are not intended for the detection task. ZeroWaste [0] is collected from a full-scale MRF, offering the most comprehensive benchmark for industrial waste detection, with over 27,000 annotated instances across 4600 images and 6000 unlabeled frames. ZeroWaste stands out for its scale, strong emphasis on detection, and structured support for fully and semi-supervised learning. We adopt it as the most suitable dataset for advancing robust, scalable object detection in real-world waste sorting.

Open-Vocabulary Object Detection. OVOD [51, 57] enables recognition of object categories not seen during training by leveraging VLMs such as CLIP [41] and ALIGN [18]. Recent approaches improve region-text alignment for open-set detection: OWL-ViT [55, 56] uses vision transformers trained on aligned image-text pairs; Grounding DINO [27] introduces language-guided detection with grounded pretraining and a strong transformer-based architecture; and YOLO-World [8] extends the YOLO family for real-time OVOD via a contrastive region-text loss and vision-language path aggregation. While these models perform well on general-purpose benchmarks like COCO [23], LVIS [14], and Objects365 [47], their effectiveness under the severe occlusion, deformation, and clutter typical of waste sorting environments remains largely untested. We address this gap by systematically benchmarking OVOD models under domain-specific industrial conditions.

Ensemble-Based Pseudo-Labeling. Pseudo-labeling (PL) reduces manual annotation costs in object detection [30, 48, 58]. Ensemble-based approaches [51, 56] further enhance pseudo-label quality by aggregating multiple predictors, mitigating errors from individual models. Most prior ensembling efforts [30, 31, 48, 58, 56] focus on clean datasets e.g., COCO [23] and VOC [11], with little exploration under dense clutter, deformation, and occlusion. We develop an ensemble-based PL pipeline for MRF waste detection, improving pseudo-label quality on the ZeroWaste-s subset and enabling scalable semi-supervised training.

Semi-Supervised Object Detection. Semi-Supervised Object Detection (SSOD) boosts detection performance by exploiting unlabeled data. Methods such as STAC [48], Unbiased Teacher [30], Soft Teacher [58], and Dense Teacher [51] use teacher-student architectures and consistency regularization. Although SSOD has shown strong results on benchmarks like COCO and VOC, it remains underexplored in dense industrial waste scenarios. We leverage ensemble-generated pseudo-labels for semi-supervised fine-tuning, establishing stronger baselines for waste recovery tasks.

3 Our Approach

This section begins with describing the ZeroWaste dataset [0] and its challenges, followed by a comprehensive zero-shot evaluation of state-of-the-art open-vocabulary detectors. We then establish new supervised baselines through fine-tuning the latest detectors on ZeroWaste-f, and finally, we introduce a semi-supervised framework that leverages ensemble-based soft pseudo-labeling to exploit the unlabeled ZeroWaste-s subset.

3.1 Dataset

All experiments are conducted on the ZeroWaste dataset, collected in an operational MRF using high-resolution (1920×1080) overhead imagery that mimics real-world industrial sort-

ing systems. Compared to conventional datasets, ZeroWaste introduces domain-specific challenges such as severe occlusions, deformations, cluttered backgrounds, and extreme class imbalance. The dataset comprises two subsets: **ZeroWaste-f**, which contains 4,503 labeled images with bounding box annotations for four recyclable categories—*cardboard*, *soft plastic*, *rigid plastic*, and *metal*; and **ZeroWaste-s**, which includes 6,212 unlabeled images captured under the same conditions to support semi-supervised learning. Fig. 1 shows examples from both subsets. A particularly challenging aspect is the class imbalance: *cardboard* accounts for over 66% of annotations, whereas *metal* comprises less than 2%. With many frames containing 15+ objects, the dataset presents high clutter and detection difficulty. These characteristics make ZeroWaste a rigorous benchmark for evaluating object detectors in unconstrained environments. While the dataset defines a fixed set of categories, it enables evaluating open-vocabulary detectors in the zero-shot setting, assessing their ability to localize and recognize unseen categories without any category-specific training.



Figure 1: Visual examples from the ZeroWaste dataset. (a) Sample images from ZeroWaste-f (top) with ground-truth annotations (bottom) for four recyclable categories: cardboard, soft plastic, rigid plastic, and metal. (b) Sample images from the unlabeled ZeroWaste-s subset, designed to support semi-supervised learning under the same real-world conditions.

To contextualize model performance, the dataset authors also reported supervised baselines using CNN-based detectors. As shown in Table 1, TridentNet [22] achieved the best performance (mAP 24.2), outperforming RetinaNet [24] and Mask R-CNN [16]. However, all models consistently scored poorly on detection metrics, highlighting the real-world difficulty of the dataset and motivating the need for updated baselines with modern architectures.

Model	mAP	mAP50	mAP75	mAPs	mAPm	mAPI
RetinaNet	21.0	33.5	22.2	4.3	9.5	22.7
Mask R-CNN	22.8	34.9	24.4	4.6	10.6	25.8
TridentNet	24.2	36.3	26.6	4.8	10.7	26.1

Table 1: Performance of CNN detectors (RetinaNet, Mask R-CNN, TridentNet) fine-tuned on ZeroWaste-f and evaluated on its test set, as reported by Bashkirova et al. [2].

3.2 Zero-Shot OVOD Baselines

To establish baseline performance for OVOD in the ZeroWaste dataset, we evaluated three state-of-the-art models: Grounding DINO [27], OWLv2 [56], and YOLO-World [8] in a zero-shot setting. These models were selected for their complementary architectures and proven capabilities. Grounding DINO achieves strong zero-shot grounding via cross-modality fusion and has set state-of-the-art benchmarks on COCO [23] and LVIS [44]; OWLv2 improves open-vocabulary performance on rare categories using large-scale self-training [56];

and YOLO-World combines vision-language modeling with real-time efficiency, making it suitable for industrial deployment [8].

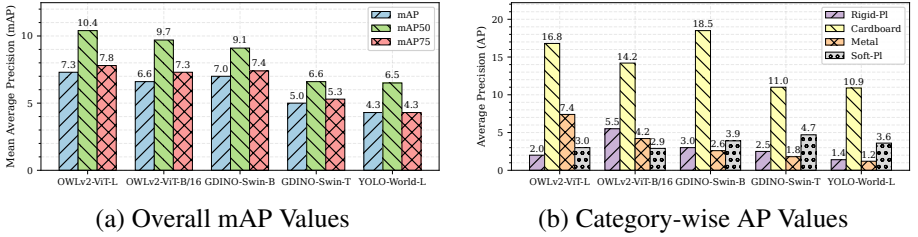


Figure 2: Zero-shot detection performance on the ZeroWaste-f test set using class-only prompts. (a) Overall mAP scores remain low across models, with OWLv2-ViT-L achieving the best mAP (7.3). (b) Category-wise AP reveals strong performance on cardboard, while soft plastic and metal suffer due to transparency and reflectivity, highlighting challenges.

Baseline Evaluation with Class-Only Prompts. To assess raw zero-shot performance, we provide each model with simple category-level prompts, namely "*cardboard*", "*soft plastic*", "*rigid plastic*", and "*metal*", which match the four predefined categories in the ZeroWaste dataset. These class-only queries contain no contextual information and serve as a minimal baseline for open-vocabulary generalization. All models are evaluated in a zero-shot setting on the ZeroWaste-f test set using COCO-style mAP metrics. Fig. 2 depicts results, revealing uniformly low performance ($\text{mAP} \leq 7.3$), with the best predictions observed for cardboard, a large and visually distinctive class. In contrast, detection of *rigid plastic*, *soft plastic* and *metal* remains particularly poor across all models, reflecting the challenge of recognizing deformable, transparent or reflective objects in cluttered industrial scenes. Further confusion-matrix analysis and prompt-level qualitative comparisons are provided in Appendix B.

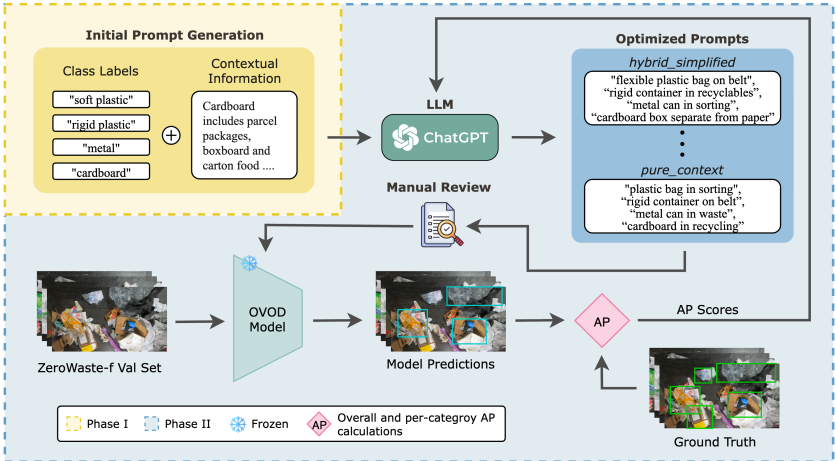


Figure 3: Iterative prompt optimization pipeline. Class-level prompts are enriched with material-specific and contextual cues using GPT-4o, then evaluated in a zero-shot setting on the ZeroWaste-f validation set using OVOD models (Grounding DINO, OWLv2). Detection performance is used to iteratively refine prompts, forming a feedback loop.

Prompt Optimization with LLMs. Fig. 3 illustrates our two-phase prompt optimization pipeline, where GPT-4o iteratively enriches class-level queries with material and contextual cues (e.g., “flexible plastic bag”). Refined prompts are evaluated on the ZeroWaste-f validation set, and the best-performing variants are used on the test set. For a detailed breakdown of prompt styles and their category-wise performance, refer to Appendix B. Fig. 5 shows that these optimized prompts markedly improve zero-shot performance for Grounding DINO (Swin-B) and OWLv2 (ViT-L), the top models from the class-only prompt evaluation. OWLv2 improves from 7.3 to 13.5 mAP (+6.2), and Grounding DINO from 7.0 to 12.4 (+5.4), with consistent gains across mAP50 and mAP75. All categories benefit, with the most significant improvements observed for *rigid plastic* and *soft plastic*. However, Grounding DINO still struggles more than OWLv2 on *metal* and *soft plastic*, indicating greater sensitivity to reflective and transparent materials. Qualitative examples illustrating the effects of prompt optimization on detection quality are in Appendix B. Despite improvements

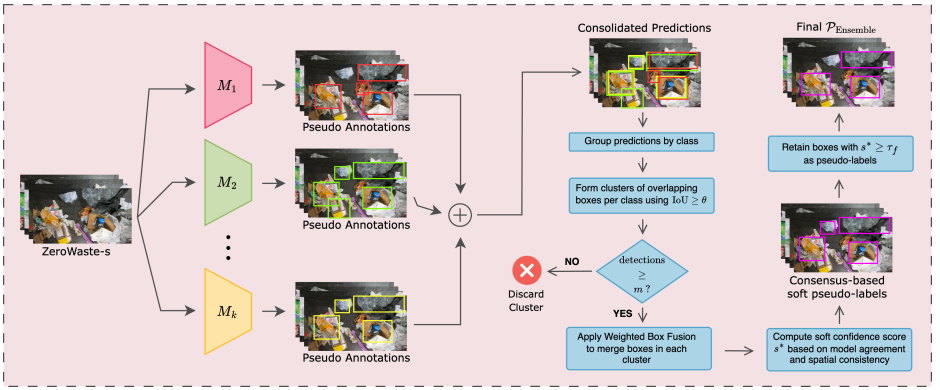


Figure 4: Soft labeling pipeline for ensemble-based pseudo-labeling. Predictions from K models are filtered, clustered by class and IoU, and fused via Weighted Box Fusion. Final confidence s^* is derived from base confidence and a consensus factor that accounts for spatial spread and model agreement. High-confidence pseudo-labels are retained for semi-supervised training.

from prompt optimization, zero-shot OVID models continue to underperform on industrial waste data, falling well below supervised baselines. Even with enriched textual prompts, these models struggle to generalize to the complex visual characteristics of industrial waste, with limitations largely driven by domain shift and insufficient exposure to waste-specific imagery during pretraining. This highlights the need for targeted task adaptation through supervised fine-tuning. For a quantitative analysis of domain shift, refer to Appendix A.

3.3 Fine-Tuning and Updated Baselines

To overcome the limitations of zero-shot OVID, we establish updated baselines through fully supervised fine-tuning on ZeroWaste-f. While the original CNN-based baselines (RetinaNet [24], Mask R-CNN [16], and TridentNet [27]) provided a solid reference point, their performance remained low under challenging MRF conditions. Motivated by recent advances in object detection, we revisit this benchmark using modern architectures. Given the fixed label space of our task, we prioritize closed-set detectors and fine-tune six state-of-the-art models: YOLO11 [10], RT-DETR [62], DINO [63], Co-DETR [69], DETA [69], and

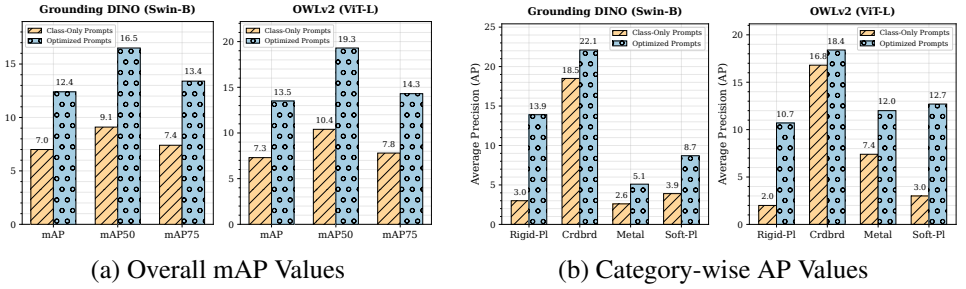


Figure 5: Class-only vs. optimized prompts on the ZeroWaste-f test set, using Grounding DINO (Swin-B) and OWLv2 (ViT-L). (a) Overall mAP, mAP50, and mAP75 show clear gains post-optimization. (b) Per-category AP indicates substantial improvements across all categories for both models.

Grounding DINO [24]. These models were selected for their strong performance on standard detection benchmarks and suitability for closed-set detection. YOLO11 and RT-DETR represent cutting-edge real-time architectures, while Grounding DINO, though originally designed for open-vocabulary detection, was included due to its strong zero-shot performance and adaptability to closed-set fine-tuning. For closed-set fine-tuning and inference, Grounding DINO used simple category-level prompts ("cardboard", "soft plastic", "rigid plastic", and "metal") exactly matching the dataset class names. The text encoder was frozen during fine-tuning, with only the visual backbone and detection head updated to adapt to waste detection. All other detectors in this comparison are purely visual and do not employ a text encoder, so they do not consume text prompts at inference.

We fine-tuned all models from their official pre-trained checkpoints on the ZeroWaste-f training split and evaluated on the test split using COCO metrics. Implementations followed Ultralytics [64] for YOLO11 and RT-DETR, MMDetection [7] for Grounding DINO, and the official repositories for DINO [63], Co-DETR [69], and DETA [69]. Except for minibatch size, we followed each codebase’s default training recipe. Training used a single NVIDIA A100 (40 GB) GPU with batch sizes of 16 (YOLO11, RT-DETR), 4 (Grounding DINO, Swin-B), and 2 (DINO, Co-DETR, DETA, all Swin-L). As shown in Table 2, fine-tuning yields substantial gains over the TridentNet baseline (24.2 mAP), with Co-DETR (Swin-L), DETA (Swin-L), and Grounding DINO (Swin-B) each achieving 51.6 mAP—more than doubling baseline performance. These results underscore the impact of task-specific fine-tuning and the superiority of transformer-based detectors over legacy CNNs and zero-shot OVO approaches. Appendix D lists the configuration identifiers, checkpoint filenames, and pre-training datasets for all fine-tuned models.

3.4 Semi-Supervised Learning with Ensemble-Based Pseudo-Labeling

While fine-tuning state-of-the-art detectors on the labeled ZeroWaste-f subset (\mathcal{D}_l) yields strong performance, scalability remains limited by the scarcity of high-quality annotations. To address this, we leverage the unlabeled ZeroWaste-s subset (\mathcal{D}_u) within a semi-supervised learning (SSL) framework. Our method constructs an ensemble $\mathcal{M} = \{M_1, M_2, \dots, M_K\}$ of fine-tuned detectors and fuses their predictions to generate soft pseudo-labels. Each pseudo-label’s confidence is refined based on spatial consistency and inter-model agreement, producing reliable supervision for \mathcal{D}_u without additional manual effort (Fig. 4). See Appendix

C for an algorithmic description and qualitative comparison with human annotations.

Model	Backbone	Venue	mAP	mAP50	mAP75	mAPs	mAPm	mAPI
TridentNet	ResNet-50	ICCV 2019	24.2	36.3	26.6	4.8	10.7	26.1
YOLO11 (L)	-	-	34.0 +9.8	45.3	36.2	7.3	16.7	36.6
RT-DETR (L)	CSPResNet-50	CVPR 2024	35.1 +10.9	45.1	37.8	7.6	16.9	37.2
DINO	ResNet-50	ICLR 2023	41.0 +16.8	51.8	44.4	9.3	19.9	44.8
DINO	Swin-L	ICLR 2023	48.3 +24.1	59.8	53.0	5.9	26.8	52.0
DETA	ResNet-50	-	38.5 +14.3	49.0	42.2	8.5	22.3	41.3
DETA	Swin-L	-	51.6 +27.4	62.6	56.1	11.6	31.5	54.6
Co-DETR	ResNet-50	ICCV 2023	37.8 +13.6	48.9	40.8	13.1	20.2	40.7
Co-DETR	Swin-L	ICCV 2023	51.6 +27.4	63.0	55.3	12.7	31.8	54.6
Grounding DINO	Swin-T	ECCV 2024	45.6 +21.4	56.8	50.1	16.6	27.3	48.9
Grounding DINO	Swin-B	ECCV 2024	51.6 +27.4	63.2	56.1	9.8	29.8	55.2

Table 2: Performance of fine-tuned object detectors on the ZeroWaste-f test set. We report mAP at standard thresholds and across object scales. TridentNet, the strongest baseline from the original ZeroWaste paper, is included for comparison. Relative mAP improvements over TridentNet are shown in green. (L) indicates large model variants.

Consensus-Based Pseudo-Label Generation. Each model $M_k \in \mathcal{M}$ predicts a set of detections on image $x \in \mathcal{D}_u$. These are first grouped by category, then clustered using pairwise IoU $\geq \theta$. For a category-specific cluster \mathcal{C} that contains detections from at least m distinct models, we apply weighted-box fusion (WBF) [49]:

$$B^* = \sum_{i=1}^{|\mathcal{C}|} w_i B_i, \quad w_i = \frac{s_i}{\sum_{j=1}^{|\mathcal{C}|} s_j} \quad (1)$$

where B_i and s_i denote the bounding box and confidence score of the i -th detection, and $|\mathcal{C}|$ is the number of detections in the cluster.

Soft Labeling via Model Agreement. We compute the soft confidence score s^* for each fused box by combining spatial consistency and model agreement:

$$\text{spread} = 1 - \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \text{IoU}(B_i, B^*), \quad (2)$$

$$\text{cf} = \exp(-\alpha \cdot \text{spread}) [1 + \beta (|\mathcal{C}_{\text{models}}| - 2)], \quad (3)$$

$$s^* = s_{\text{base}} \cdot \text{cf} \quad (4)$$

Here, $|\mathcal{C}_{\text{models}}|$ is the number of contributing models, and $s_{\text{base}} = \max_{i \in \mathcal{C}} s_i$.

For hard pseudo-labels, s^* is set as s_{base} , without adjusting for spatial consistency or model agreement. Clusters with $s^* \geq \tau_f$ form the pseudo-label set $\mathcal{P}_{\text{Ensemble}}$. Parameters α and β penalize spatial spread and reward model consensus, suppressing noisy predictions.

Semi-Supervised Training. A detector \mathcal{F} is trained on mixed mini-batches of \mathcal{D}_l and pseudo-labeled \mathcal{D}_u (ratio 1:2). Pseudo-labeled samples are weighted by $(s^*)^p$ and the total training loss becomes:

$$\mathcal{L}_{\text{unsup}} = \sum_i (s_i^*)^p [\mathcal{L}_{\text{cls}} + \lambda \mathcal{L}_{\text{reg}}] \quad (5)$$

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{sup}} + \gamma \mathcal{L}_{\text{unsup}} \quad (6)$$

Implementation Details. The ensemble \mathcal{M} comprises Grounding DINO (Swin-B), Co-DETR (Swin-L), DETA (Swin-L), and DINO (Swin-L), fine-tuned on \mathcal{D}_l . Predictions are filtered with $\tau_f = 0.35$, clustered with $\theta = 0.65$, and retained if $m \geq 2$ model agreement is observed. Remaining hyper-parameters are $\alpha = 5.0$, $\beta = 0.1$, $p = 2.0$, $\lambda = 2.0$, and $\gamma = 1.0$. All experiments were conducted using Grounding DINO (Swin-B), which achieved the overall highest performance in the fully supervised setting (Table 2).

Results. Table 3 summarizes the effect of different pseudo-labeling strategies. For both Swin-T and Swin-B backbones, hard pseudo-labels $\mathcal{P}_{\text{Swin-B}}$, generated by a single fine-tuned Grounding DINO (Swin-B) model, yield modest improvements over the fully supervised baseline (+2.1 and +0.9 mAP, respectively). Using ensemble-generated hard pseudo-labels further improves performance, but the best results are achieved with our proposed soft ensemble labeling, which boosts mAP to 49.3 for Swin-T and 54.3 for Swin-B, representing +3.7 and +2.7 mAP gains, respectively. To assess whether these gains disproportionately favor majority classes, Table 4 reports per-class AP (AP@[50:95]): our soft ensemble improves all classes for both backbones, including the rare *metal* class (+1.5 AP on Swin-T, +1.9 on Swin-B), with the largest gains observed on rigid and soft plastic.

Training Type	Label Type	Pseudo-Label Source	Backbone	ZeroWaste-f mAP			Improvement over Supervised		
				mAP	mAP@50	mAP@75	Δ mAP	Δ @50	Δ @75
Fully Supervised	Ground Truth	—	Swin-T	45.6	56.8	50.1	—	—	—
Semi-Supervised	Hard	$\mathcal{P}_{\text{Swin-B}}$	Swin-T	47.7	58.1	52.3	+2.1	+1.3	+2.2
Semi-Supervised	Hard	$\mathcal{P}_{\text{Ensemble (Ours)}}$	Swin-T	48.6	59.7	53.2	+3.0	+2.9	+3.1
Semi-Supervised	Soft	$\mathcal{P}_{\text{Ensemble (Ours)}}$	Swin-T	49.3	60.3	54.1	+3.7	+3.5	+4.0
Fully Supervised	Ground Truth	—	Swin-B	51.6	63.2	56.0	—	—	—
Semi-Supervised	Hard	$\mathcal{P}_{\text{Swin-B}}$	Swin-B	52.5	64.2	57.3	+0.9	+1.0	+1.3
Semi-Supervised	Hard	$\mathcal{P}_{\text{Ensemble (Ours)}}$	Swin-B	53.8	65.5	59.0	+2.2	+2.3	+3.0
Semi-Supervised	Soft	$\mathcal{P}_{\text{Ensemble (Ours)}}$	Swin-B	54.3	65.9	59.5	+2.7	+2.7	+3.5

Table 3: Comparison of supervised and semi-supervised training on the ZeroWaste-f test set using Grounding DINO with Swin-T and Swin-B backbones. Hard pseudo-labels from a single model ($\mathcal{P}_{\text{Swin-B}}$) offer moderate gains, while our ensemble-based soft labels yield the highest performance. Relative improvements over full supervision are shown in the right-most columns.

Training Type	Label Type	Pseudo-Label Source	Backbone	ZeroWaste-f AP (per class)			
				Rigid Plastic	Cardboard	Metal	Soft Plastic
Fully Supervised	Ground Truth	—	Swin-T	48.9	50.2	36.6	47.8
Semi-Supervised	Soft	$\mathcal{P}_{\text{Ensemble (Ours)}}$	Swin-T	53.3 ^{+4.4}	52.9 ^{+2.7}	38.1 ^{+1.5}	52.9 ^{+5.1}
Fully Supervised	Ground Truth	—	Swin-B	57.0	54.6	42.1	53.0
Semi-Supervised	Soft	$\mathcal{P}_{\text{Ensemble (Ours)}}$	Swin-B	61.2 ^{+4.2}	55.8 ^{+1.2}	44.0 ^{+1.9}	56.1 ^{+3.1}

Table 4: Per-class AP on ZeroWaste-f test set for Grounding DINO under fully supervised vs. semi-supervised (soft, ensemble) training. Semi-supervised (soft, ensemble) improves performance across all classes, including rare ones such as *metal* and *soft plastic*.

4 Final Pseudo-Annotations for ZeroWaste-s

Having demonstrated the effectiveness of our ensemble-based soft pseudo-labeling (PL) strategy, we use the best-performing model from Table 3, Grounding DINO (Swin-B), trained with soft $\mathcal{P}_{\text{Ensemble}}$ pseudo-labels, to generate a final set of high-quality annotations for the

unlabeled ZeroWaste-s subset. The resulting pseudo-labels, denoted as $\mathcal{P}_{\text{Final}}$, are intended to support research in semi-supervised object detection and enhance the utility of ZeroWaste-s as a high-quality dataset for real-world industrial waste sorting applications, helping address the broader scarcity of annotated datasets in this domain. To ensure label reliability, we apply a conservative confidence threshold of 0.4, balancing precision and category diversity. The final $\mathcal{P}_{\text{Final}}$ comprises 33,075 bounding boxes across 6,065 images, with category-level statistics in Table 5.

Category	Annotations
Cardboard	21,352
Soft Plastic	9,806
Rigid Plastic	1,523
Metal	394
Total	33,075

Table 5: Distribution of final pseudo-annotations.

Model	Training	mAP	mAP@50	mAP@75
YOLO11 (Large)	ZeroWaste-f	34.0	45.3	36.2
YOLO11 (Large)	$\mathcal{P}_{\text{Final}}$	40.3 +6.3	51.1 +5.8	43.3 +7.1
RT-DETR (Large)	ZeroWaste-f	35.1	45.1	37.8
RT-DETR (Large)	$\mathcal{P}_{\text{Final}}$	39.4 +4.3	52.1 +7.0	42.6 +4.8

Table 6: YOLO11 and RT-DETR trained on labeled (ZeroWaste-f) vs. pseudo-labeled ($\mathcal{P}_{\text{Final}}$) data, evaluated on ZeroWaste-f test set.

Indirect Evaluation via Model Transfer. To evaluate the generalization quality of $\mathcal{P}_{\text{Final}}$, we fine-tune two high-performing detectors (YOLO11 and RT-DETR) exclusively on this pseudo-labeled dataset and evaluate their performance on the labeled ZeroWaste-f test set. This setup simulates a practical use case where manually labeled data is scarce, and model training depends entirely on pseudo-annotations. Table 6 compares the performance of detectors trained on $\mathcal{P}_{\text{Final}}$ versus those trained on manually labeled ZeroWaste-f data. The results confirm that $\mathcal{P}_{\text{Final}}$ can serve as an effective alternative to manual annotation. YOLO11 trained on pseudo-labeled data achieves +6.3 mAP improvement over its fully supervised baseline, while RT-DETR sees a +4.3 mAP gain. The consistent improvements across mAP, mAP@50, and mAP@75 indicate not only better detection accuracy but also enhanced localization precision. These findings validate the generalization strength and practical utility of our pseudo-annotation pipeline in real-world applications.

5 Conclusion

We present a comprehensive framework for advancing AI-driven waste detection in industrial settings by evaluating zero-shot open-vocabulary object detectors, fine-tuning state-of-the-art models, and introducing a robust semi-supervised learning pipeline using ensemble-based soft pseudo-labeling. Our extensive experiments on the challenging ZeroWaste dataset reveal that while current OVOD models struggle in cluttered, deformable environments, targeted prompt optimization and fine-tuning yield substantial performance gains. Furthermore, our consensus-driven pseudo-labeling approach enables scalable learning from unlabeled data and produces high-quality annotations that rival or surpass fully supervised baselines when used for training. These contributions establish new benchmarks and outline a scalable path toward AI-assisted waste recovery in real-world material recovery facilities.

References

- [1] AMP Robotics. AMP Robotics. <https://www.amrobotics.com/>. Accessed: 2020-05-30.
- [2] Dina Bashkirova, Olga Russakovsky, and Stella X. Yu. Zerowaste dataset: Towards deformable object segmentation in cluttered scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12725–12734, 2022. URL https://openaccess.thecvf.com/content/CVPR2022/papers/Bashkirova_ZeroWaste_Dataset_Towards_Deformable_Object_Segmentation_in_Cluttered_Scenes_CVPR_2022_paper.pdf.
- [3] C. Bircanoglu, M. Atay, F. Beser, O. Genc, and M. A. Kizrak. Recyclenet: Intelligent waste sorting using deep neural networks. In *Proceedings of the International Conference on Innovative Intelligent Systems and Applications (INISTA)*, pages 1–7, Jul. 2018. doi: 10.1109/INISTA.2018.8466276.
- [4] Encyclopædia Britannica. Materials recovery facility. Online, n.d. URL <https://www.britannica.com/technology/materials-recovery-facility>. Accessed: 20-02-2025.
- [5] B. D. Carolis, F. Ladogana, and N. Macchiarulo. Yolo trashnet: Garbage detection in video streams. In *Proceedings of the IEEE Conference on Evolving and Adaptive Intelligent Systems (EAIS)*, pages 1–7, May 2020. doi: 10.1109/EAIS48028.2020.9122693.
- [6] Sara Casao, Fernando Peña, Alberto Sabater, Rosa Castellón, Darío Suárez, Eduardo Montijano, and Ana C. Murillo. Spectralwaste dataset: Multimodal data for waste sorting automation. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5852–5858, 2024. doi: 10.1109/IROS58592.2024.10801797.
- [7] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- [8] Tianheng Cheng, Lin Song, Yixiao Ge, Wenyu Liu, Xinggang Wang, and Ying Shan. Yolo-world: Real-time open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16901–16911, June 2024.
- [9] Datacluster-labs. Domestic trash dataset. <https://github.com/datacluster-labs/Domestic-Trash-Dataset>, 2021. Accessed: 2025-03-01.
- [10] Environmental Protection Agency (EPA). Advancing sustainable materials management: 2018 fact sheet, 2020. URL <https://www.epa.gov/facts-and-figures-about-materials-waste-and-recycling/national-overview-facts-and-figures-materials>.

- [11] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2): 303–338, June 2010.
- [12] Zhicheng Feng, Jie Yang, Lifang Chen, Zhichao Chen, and Linhong Li. An intelligent waste-sorting and recycling device based on improved efficientnet. *International Journal of Environmental Research and Public Health*, 19(23):15987, 2022. doi: 10.3390/ijerph192315987.
- [13] Sathish Paulraj Gundupalli, Subrata Hait, and Atul Thakur. A review on automated sorting of source-separated municipal solid waste for recycling. *Waste Management*, 60:56–74, 2017. doi: 10.1016/j.wasman.2016.09.015.
- [14] Agrim Gupta, Piotr Dollár, Ross Girshick, et al. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [15] M. Haamer. Wade-ai dataset. Available at Wade-ai project page, 2020.
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [17] Nidhal Jegham, Chan Young Koh, Marwan Abdelatti, and Abdeltawab Hendawi. Yolo evolution: A comprehensive benchmark and architectural review of yolov12, yolov11, and their previous versions, 2025. URL <https://arxiv.org/abs/2411.00201>.
- [18] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision, 2021. URL <https://arxiv.org/abs/2102.05918>.
- [19] Silpa Kaza, Lisa Yao, Perinaz Bhada-Tata, and Frank Van Woerden. *What a Waste 2.0: A Global Snapshot of Solid Waste Management to 2050*. World Bank, 2018. URL <https://datatopics.worldbank.org/what-a-waste/>. Accessed: 2020-05-30.
- [20] Radka Kopecká, Marlies Hrad, and Marion Huber-Humer. The role of the waste sector in the sustainable development goals and the ipcc assessment reports. *Österreichische Wasser- und Abfallwirtschaft*, 76:300–307, 2024. doi: 10.1007/s00506-024-01034-7.
- [21] Maria Koskinopoulou, Fredy Raptopoulos, George Papadopoulos, Nikitas Mavrakis, and Michail Maniadakis. Robotic waste sorting technology: Toward a vision-based categorization system for the industrial robotic separation of recyclable waste. *IEEE Robotics & Automation Magazine*, 28(2):50–60, 2021.
- [22] Yanghao Li, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Scale-aware trident networks for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6054–6063, 2019.

- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. *European Conference on Computer Vision (ECCV)*, 2014.
- [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [25] W. Lin. Yolo-green: A real-time classification and object detection model optimized for waste management. In *Proceedings of the IEEE International Conference on Big Data (Big Data)*, pages 51–57, Dec. 2021. doi: 10.1109/BIGDATA52589.2021.9671821.
- [26] F. Liu, H. Xu, M. Qi, D. Liu, J. Wang, and J. Kong. Depth-wise separable convolution attention module for garbage image classification. *Sustainability*, 14(5):3099, Mar. 2022. doi: 10.3390/su14053099.
- [27] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision*, pages 38–55. Springer, 2024.
- [28] W. Liu, H. Ouyang, Q. Liu, S. Cai, C. Wang, J. Xie, and W. Hu. Image recognition for garbage classification based on transfer learning and model fusion. *Mathematical Problems in Engineering*, 2022:1–12, Aug. 2022. doi: 10.1155/2022/4793555.
- [29] Y. Liu, Z. Ge, G. Lv, and S. Wang. Research on automatic garbage detection system based on deep learning and narrowband internet of things. *Journal of Physics: Conference Series*, 1069, Aug. 2018. doi: 10.1088/1742-6596/1069/1/012032.
- [30] Yen-Cheng Liu, Chia-Yi Ma, Zijian He, Simon Kuo, and Jia-Bin Huang. Unbiased teacher for semi-supervised object detection. In *International Conference on Learning Representations (ICLR)*, 2021.
- [31] Yue Liu, Shunping Wang, et al. Dense teacher: Dense pseudo-labels for semi-supervised object detection. In *International Conference on Learning Representations (ICLR)*, 2022.
- [32] W. Ma, X. Wang, and J. Yu. A lightweight feature fusion single shot multibox detector for garbage detection. *IEEE Access*, 8:188577–188586, 2020. doi: 10.1109/ACCESS.2020.3031990.
- [33] Anthony Martin. Recycling image classification. Online. URL <http://web.cecs.pdx.edu/~singh/rcyc-web/index.html>. Accessed: 20-02-2025.
- [34] O. A. Mengistu. Smart trash net: Waste localization and classification. Online, 2017. URL <https://www.semanticscholar.org/paper/Final-Report-%3A-Smart-Trash-Net-%3A-Waste-LocalizationAwe-Mengistu/581fb0f0405c7f0e60610d88ceaceb9af44d8569>. Accessed: 22-02-2025.
- [35] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. In *European conference on computer vision*, pages 728–755. Springer, 2022.

- [36] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection. *Advances in Neural Information Processing Systems*, 36:72983–73007, 2023.
- [37] G. Mittal, K. B. Yagnik, M. Garg, and N. C. Krishnan. Garbage in images (gini) dataset. Online, 2016. URL <https://github.com/spotgarbage/spotgarbage-GINI>. Accessed: 22-02-2025.
- [38] United Nations. Sustainable development goal 12: Responsible consumption and production. Online, n.d. URL <https://sdgs.un.org/goals/goal12>. Accessed: 19-02-2025.
- [39] Jeffrey Ouyang-Zhang, Jang Hyun Cho, Xingyi Zhou, and Philipp Krähenbühl. Nms strikes back, 2022. URL <https://arxiv.org/abs/2212.06137>.
- [40] P. F. Proença and P. Simões. Taco: Trash annotations in context for litter detection. 2020. URL <https://arxiv.org/abs/2003.06975>.
- [41] Alec Radford, Jong Wook Kim, Christopher Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021.
- [42] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection, 2016. URL <https://arxiv.org/abs/1506.02640>.
- [43] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks, 2016. URL <https://arxiv.org/abs/1506.01497>.
- [44] AMP Robotics. Ai-powered robotics for recycling. *Recycling News*, 2022.
- [45] Aleksei Seredkin et al. Automated waste sorting: object detection using convolutional neural networks. *IOP Conference Series: Earth and Environmental Science*, 337(1): 012048, 2019.
- [46] A. Serezhkin. Drinking waste classification dataset. <https://www.kaggle.com/datasets/arkadiyhacks/drinking-waste-classification>, 2020.
- [47] Shuai Shao, Zeming Zhao, Bo Li, et al. Objects365: A large-scale, high-quality dataset for object detection. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [48] Kihyuk Sohn, Zizhao Zhang, Chun-Liang Li, Han Zhang, Chen-Yu Lee, and Tomas Pfister. A simple semi-supervised learning framework for object detection. *arXiv preprint*, 2020. URL <https://arxiv.org/abs/2005.04757>.
- [49] Roman Solovyev, Weimin Wang, and Tatiana Gabruseva. Weighted boxes fusion: Ensembling boxes from different object detection models. *Image and Vision Computing*, 107:104117, March 2021. ISSN 0262-8856. doi: 10.1016/j.imavis.2021.104117. URL <http://dx.doi.org/10.1016/j.imavis.2021.104117>.

- [50] Joao Sousa, Ana Rebelo, and Jaime S. Cardoso. Automation of waste sorting with deep learning. In *2019 XV Workshop de Visão Computacional (WVC)*, pages 43–48. IEEE, 2019.
- [51] A. Sun and H. Xiao. Thanosnet: A novel trash classification method using metadata. In *Proceedings of the IEEE International Conference on Big Data (Big Data)*, pages 1394–1401, Dec. 2020. doi: 10.1109/BigData50022.2020.9378287.
- [52] Xueyong Tian, Liwei Shi, Yuanqing Luo, and Xinlong Zhang. Garbage classification algorithm based on improved mobilenetv3. *IEEE Access*, 2023.
- [53] Maria Triassi, Rita Alfano, Maddalena Illario, Antonio Nardone, Ornella Caporale, and Paolo Montuori. Environmental pollution from illegal waste disposal and health effects: A review on the “triangle of death”. *International Journal of Environmental Research and Public Health*, 12:1216–1236, 2015. doi: 10.3390/ijerph120201216.
- [54] Ultralytics. Models supported by ultralytics yolo. <https://docs.ultralytics.com/models/>, 2025. Accessed: 2025-08-15.
- [55] A. H. Vo, L. H. Son, M. T. Vo, and T. Le. A novel framework for trash classification using deep transfer learning. *IEEE Access*, 7:178631–178639, 2019. doi: 10.1109/ACCESS.2019.2959033.
- [56] Waste-Robotics. Waste-Robotics. <https://wasterobotic.com/>. Accessed: 2020-05-30.
- [57] Zhongyi Xia, Houkui Zhou, Huimin Hu, Haoji Hu, Guangqun Zhang, Junguo Hu, and Tao He. YOLO-MTG: A lightweight YOLO model for multi-target garbage detection. *Signal, Image and Video Processing*, pages 1–16, 2024. doi: 10.1007/s11760-024-0320-2. URL <https://link.springer.com/article/10.1007/s11760-024-03220-2>.
- [58] Bowen Xu, Hang Shi, Yutong Wu, Chenyang Li, Zizhao Zhang, et al. End-to-end semi-supervised object detection with soft teacher. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [59] M. Yang and G. Thung. Classification of trash for recyclability status. Technical Report 3, CS229 Project Report, Stanford University, 2016. URL <http://cs229.stanford.edu/proj2016/report/ThungYang-ClassificationOfTrashForRecyclabilityStatus-report.pdf>. Accessed: 20-02-2025.
- [60] Dmitry Yudin, Nikita Zakharenko, Artem Smetanin, Roman Filonov, Margarita Kichik, Vladislav Kuznetsov, Dmitry Larichev, Evgeny Gudov, Semen Budenny, and Aleksandr Panov. Hierarchical waste detection with weakly supervised segmentation in images from recycling plants. *Engineering Applications of Artificial Intelligence*, 128: 107542, 2024.
- [61] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14393–14402, 2021. URL https://openaccess.thecvf.com/content/CVPR2021/html/Zareian_Open-Vocabulary_Object_Detection_Using_Captions_CVPR_2021_paper.html.

- [62] Zen Robotics. Zen Robotics. <https://zenrobotics.com/>. Accessed: 2020-05-30.
- [63] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel Ni, and Heung-Yeung Shum. DINO: DETR with improved denoising anchor boxes for end-to-end object detection. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=3mRwyG5one>.
- [64] Yian Zhao, Wenyu Lv, Shangliang Xu, Jinman Wei, Guanzhong Wang, Qingqing Dang, Yi Liu, and Jie Chen. Detrs beat yolos on real-time object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16965–16974, June 2024.
- [65] P. Zhou, Z. Zhu, X. Xu, X. Liu, B. He, and J. Zhang. Towards the urban future: A novel trash segregation algorithm based on improved yolov4. In *Proceedings of the IEEE International Conference on Robotics and Biomimetics (ROBIO)*, pages 1526–1531, Dec. 2021. doi: 10.1109/ROBIO54168.2021.9739288.
- [66] Qiang Zhou, Chaohui Yu, Zhibin Wang, Qi Qian, and Hao Li. Instant-teaching: An end-to-end semi-supervised object detection framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4081–4090, 2021. URL https://openaccess.thecvf.com/content/CVPR2021/html/Zhou_Instant-Teaching_An_End-to-End_Semi-Supervised_Object_Detection_Framework_CVPR_2021_paper.html.
- [67] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision, 2022. URL <https://arxiv.org/abs/2201.02605>.
- [68] Zhen Zhou, Xiaofeng Jin, Lu Chen, and Yajuan Han. Construction waste object detection based on improved yolov5 algorithm. *Sensors*, 23(4):1987, 2023.
- [69] Zhuofan Zong, Guanglu Song, and Yu Liu. Detrs with collaborative hybrid assignments training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6748–6758, October 2023.

Appendix

We provide extended analysis to support our main findings across four sections. **Appendix A** quantifies cross-dataset domain shift between natural and industrial imagery using Maximum Mean Discrepancy (MMD) and t-SNE visualizations. **Appendix B** offers a detailed breakdown of zero-shot Open-Vocabulary Object Detection (OVOD) performance, including confusion matrices, prompt optimization, and qualitative comparisons. **Appendix C** presents our ensemble-based pseudo-labeling algorithm in full and visualizes its alignment with human annotations. **Appendix D** lists the configuration identifiers, checkpoint filenames, and pre-training datasets used to initialize each fine-tuned detector to ensure reproducibility.

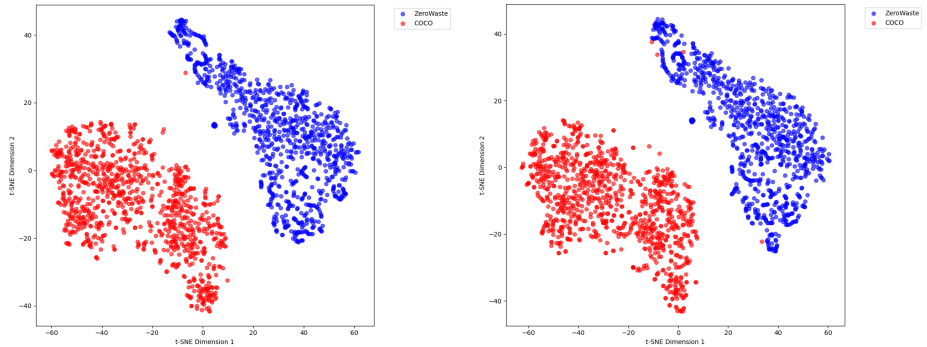
A Cross-Dataset Domain Shift: COCO vs. ZeroWaste

To validate our hypothesis that zero-shot OVOD models underperform on ZeroWaste due to domain shift, we quantify the cross-dataset distributional divergence using MMD. These models are typically pre-trained on large-scale, general-purpose natural image datasets with image-text alignment, such as COCO, OpenImages, or LAION. However, their performance may degrade when applied to specialized domains like industrial waste, which differ significantly in visual characteristics. MMD provides a kernel-based statistical measure for comparing feature distributions in a reproducing kernel Hilbert space (RKHS). To assess distributional shift, we extract deep feature embeddings from a YOLO model pre-trained on COCO and compute MMD between random samples from COCO (train or val) and the ZeroWaste-f test set.

Run	COCO (Train) vs. ZeroWaste-f (Test)	COCO (Val) vs. ZeroWaste-f (Test)
Run 1	0.6125	0.6054
Run 2	0.6099	0.6068
Run 3	0.6082	0.6077
Run 4	0.5974	0.6046
Run 5	0.6111	0.6004
Mean \pm Std	0.6078 \pm 0.0054	0.6050 \pm 0.0026

Table 7: MMD scores across five runs comparing COCO (train and val) with the ZeroWaste-f test set. The consistently high inter-dataset MMD confirms a substantial domain shift between natural and industrial image domains.

The consistently high MMD scores (~ 0.61) across both COCO train and validation splits quantitatively confirm a substantial domain gap between natural and industrial imagery. Unlike COCO’s diverse, object-centric scenes, ZeroWaste contains cluttered layouts, deformable objects, and domain-specific materials, factors that hinder generalization in zero-shot OVOD models pre-trained on natural image distributions. To visualize this shift, we apply t-SNE to the deep embeddings from both datasets. As shown in Figure 6, COCO and ZeroWaste embeddings form well-separated clusters with minimal overlap, further substantiating the observed cross-domain discrepancy.



(a) COCO (Train) vs. ZeroWaste-f (Test)

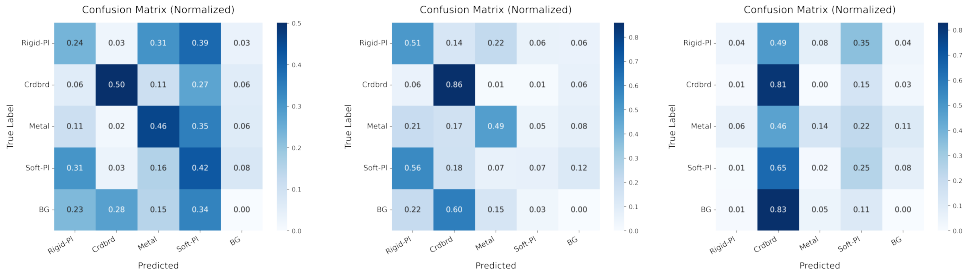
(b) COCO (Val) vs. ZeroWaste-f (Test)

Figure 6: t-SNE visualizations of deep feature embeddings extracted from a YOLO model pre-trained on COCO. Red points correspond to COCO; blue points to ZeroWaste-f. The clear cluster separation in both (a) and (b) reflects a pronounced domain shift, corroborating the high MMD scores and explaining the performance drop in zero-shot OVOD settings.

B Extended Analysis of Zero-Shot OVOD Performance

Confusion Matrix Analysis of Zero-Shot OVOD Performance. Figure 7 shows confusion matrices for Grounding DINO (Swin-B), OWLv2 (ViT-L), and YOLO-World-L on the ZeroWaste-f test set in a zero-shot setting. Cardboard is the most accurately detected class across models, with OWLv2 reaching 86% accuracy. In contrast, soft plastic, rigid plastic, and cardboard are frequently confused due to their flexible, translucent, and visually similar appearances in cluttered scenes. YOLO-World shows the highest confusion, misclassifying 65% of soft plastic as cardboard. Metal detection is the weakest across models, with over 50% misclassification, likely due to deformation and reflective surfaces. Additionally, background regions—often containing paper and conveyor belt textures resembling cardboard—are mislabeled at high rates (60–83%), especially by OWLv2 and YOLO-World. These trends highlight the challenges current OVOD models face in fine-grained categorization and generalization under real-world domain shift.

Prompt Optimization Analysis. Tables 8 and 9 present a detailed breakdown of textual query styles evaluated on the ZeroWaste-f validation set for OWLv2 (ViT-L) and Grounding DINO (Swin-B). While class-only queries offer a minimal baseline (10.2 mAP), incorporating structured, context-rich descriptions significantly improves detection. The best-performing style, *combined_success*, boosts mAP to 13.3 for OWLv2 and 13.4 for Grounding DINO, demonstrating consistent gains across categories, particularly for soft and rigid plastics. Improvements stem from precise material cues (“rigid plastic container”), contextual grounding (“on conveyor belt”), and semantic redundancy (“plastic bag or wrap”). However, overly verbose or visibility-focused prompts sometimes degrade performance, especially in OWLv2, highlighting model-specific sensitivities to linguistic style. Notably, OWLv2 benefits more from spatial descriptors, while Grounding DINO responds better to concise object-centric queries. These findings emphasize the importance of prompt engineering in adapting OVOD models to complex, domain-specific settings like industrial waste sorting.



(a) Grounding DINO (Swin-B) (b) OWLv2 (ViT-L) (c) YOLO-World-L

Figure 7: Confusion matrices of zero-shot OVID models on the ZeroWaste-f test set. (a) Grounding DINO shows confusion between cardboard and rigid plastic. (b) OWLv2 performs well on cardboard but struggles with metal. (c) YOLO-World-L shows greater overlap between soft plastic and rigid plastic classes.

Prompt Style	Description	Optimized Prompt	mAP	mAP50	AP per Category			
					Rigid-Pl	Crdbrd	Metal	Soft-Pl
<i>class-only</i>	(Baseline) Exact category labels	"soft plastic" "rigid plastic" "metal" "cardboard"	10.2	15.5	2.5	16	21.4	1.0
<i>enhanced_properties</i>	Detailed material properties with location	"flexible translucent bag on belt" "rigid hollow container in sorting" "metallic can in stream" "structured cardboard in waste"	6.9	11.8	3.7	9.0	8.2	6.8
<i>visibility_focused</i>	Combines visibility cues with context	"flexible plastic visible in sorting" "rigid container visible on belt" "metal can visible in stream" "cardboard visible in waste"	8.5	13.5	4.3	6.4	22.4	0.6
<i>pure_context</i>	Minimal descriptors with strong context	"plastic bag in sorting" "rigid container on belt" "metal can in waste" "cardboard in recycling"	8.7	14.4	5.1	12.4	14.6	2.7
<i>recycling_context</i>	Recycling-specific context	"recyclable plastic bag" "recyclable plastic container" "recyclable metal can" "recyclable cardboard box"	9.1	12.6	3.7	20.8	6.7	5.0
<i>hybrid_simplified</i>	Combines successful elements with simplification	"flexible plastic bag on belt" "rigid container in recyclables" "metal can in sorting" "cardboard box separate from paper"	9.2	15.8	7.8	17.9	10.8	0.4
<i>location_enhanced</i>	Strong emphasis on location and sorting context	"plastic bag on sorting belt" "rigid container in recycling line" "metal can in waste stream" "cardboard on conveyor belt"	10.2	15.9	7.3	9.3	23.1	1.1
<i>combined_success</i>	Combines elements from the most successful trials	"translucent plastic bag or film in waste stream" "sturdy rigid plastic container or jug in facility sorting" "shiny metal tin or can on the waste belt" "brown cardboard box separated from paper"	13.3 +3.1	20.5 +5.0	2.9	17.1	24.6	8.4

Table 8: Zero-shot evaluation results of OWLv2 (ViT-L) on the ZeroWaste-f validation set using different textual query styles. mAP across multiple IoU thresholds and per-category AP for rigid plastic (Rigid-Pl), cardboard (Crdbrd), metal (Metal), and soft plastic (Soft-Pl) are reported. The table demonstrates the impact of textual prompt refinement on model performance, highlighting the effectiveness of various prompt styles. The best-performing prompt style (*combined_success*) shows a relative improvement of +3.1 mAP over the baseline (class-only query). Reported improvements (e.g., +3.1) indicate absolute gains in overall mAP compared to the class-only prompt baseline.

Prompt Style	Description	Optimized Prompts	mAP	mAP50	AP per Category			
					Rigid-Pl	Crdbrd	Metal	Soft-Pl
class-only	(Baseline) Exact category labels	"soft plastic" "rigid plastic" "metal" "cardboard"	10.2	12.8	5.2	17.9	15.0	2.8
location_enhanced	Strong emphasis on location and sorting context	"plastic bag among waste" "plastic container in clutter" "metal can on conveyor" "cardboard box in sorting"	9.0	12.6	5.4	17.8	9.8	3.2
pure_context	Minimal descriptors with strong context	"plastic bag among waste" "hard solid plastic container" "metal can" "brown cardboard box"	9.0	12.8	4.4	21.1	5.4	5.1
recycling_context	Recycling-specific context	"recyclable plastic bag" "recyclable plastic container" "recyclable metal can" "recyclable cardboard box"	9.1	12.6	3.7	20.8	6.7	5.0
hybrid_simplified	Combines successful elements with simplification	"thin plastic bag" "hard plastic container" "metal can" "brown cardboard box"	9.1	12.8	3.2	19.0	8.6	5.6
enhanced_properties	Detailed material properties with location	"crumpled translucent plastic bag or wrap" "solid rigid plastic container or bottle" "shiny reflective metal can or tin" "thick brown cardboard box or packaging"	10.2	13.7	6.7	17.9	10.9	5.2
visibility_focused	Combines visibility cues with context	"transparent wrinkled plastic" "hard shaped plastic" "metallic shiny object" "brown flat cardboard"	11.2	15.1	4.7	20.3	14.1	5.7
combined_success	Combines elements from the most successful trials	"flexible plastic bag or wrap" "hollow rigid plastic container or bottle" "shiny metallic can" "stiff brown cardboard box"	13.4 +3.2	18.1 +5.3	8.8	22.1	16.4	6.3

Table 9: Zero-shot evaluation of Grounding DINO (Swin-B) with different text prompt styles under the same settings as Table 8. Similarly, the table highlights the impact of various prompt styles on model performance, demonstrating that the best-performing prompt style (*combined_success*) achieved a relative improvement of +3.2 mAP over the baseline (class-only prompt), as indicated in green. All reported gains reflect improvements in mAP relative to the class-only baseline.

Qualitative Results for Zero-Shot OVOE Evaluation. Figure 8 presents a visual comparison of model predictions using class-only versus optimized prompts on the ZeroWaste-f dataset. Ground truth annotations are shown on the left, while predictions from class-only and optimized prompts appear in the center and right, respectively. For both OWLv2 (ViT-L) and Grounding DINO (Swin-B), optimized prompts result in improved localization and reduced false positives, particularly for deformable classes such as soft plastic and cardboard. OWLv2 shows notable gains in recall and accuracy with enhanced prompts, while Grounding DINO exhibits reduced over-detection and cleaner bounding boxes. Nonetheless, both models consistently fail to detect metallic objects, and occluded items remain challenging. These examples illustrate the practical benefits and limitations of prompt optimization for open-vocabulary detection in cluttered, real-world waste environments.

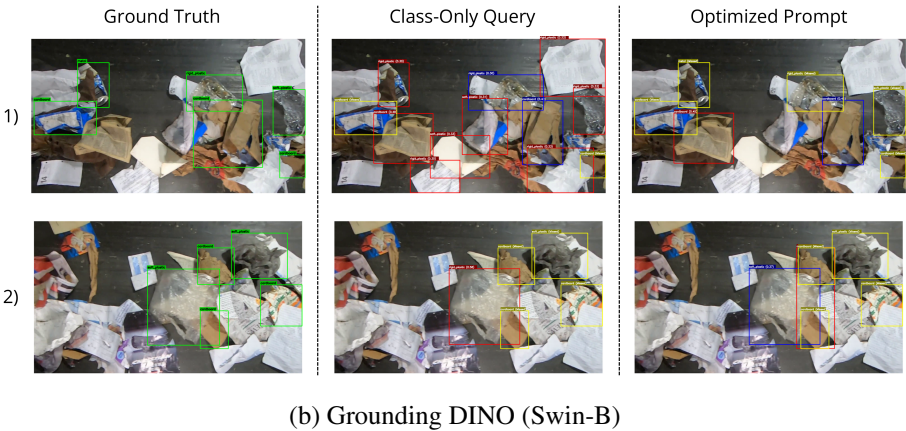
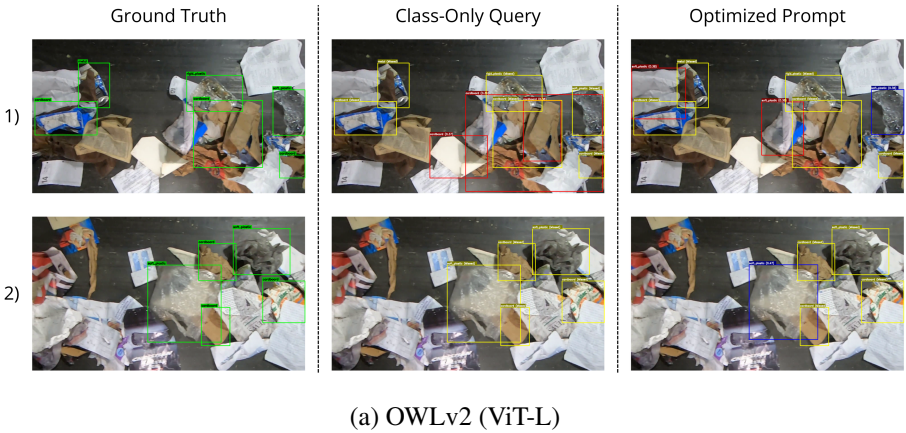


Figure 8: Qualitative comparison of object detection using class-only and optimized prompts on the ZeroWaste-f dataset. Ground truth annotations (green) are shown on the left, with model predictions using class-only (middle) and optimized prompts (right). Correct detections are in blue, incorrect in red, and missed detections in yellow. Optimized prompts improve localization, particularly for soft plastic and cardboard, while reducing false positives. However, both models struggle with metal detection and occluded objects. Results are shown for (a) OWLv2 (ViT-L) and (b) Grounding DINO (Swin-B).

C Ensemble-Based Pseudo-Labeling: Algorithm and Qualitative Analysis

Algorithmic Details of Ensemble-Based Pseudo-Labeling. To complement the description in Section 3.4 of the main paper, Algorithm 1 formally outlines our ensemble-based soft pseudo-labeling pipeline.

Algorithm 1 Ensemble-Based Soft Pseudo-Labeling

Require: \mathcal{M} : Set of models $\{M_1, M_2, \dots, M_k\}$, \mathcal{D} : Unlabeled dataset

Require: τ : Initial confidence threshold, θ : IoU threshold, m : Minimum model agreement

Require: τ_f : Final soft confidence threshold, α : spread decay factor, β : model agreement bonus factor

```

1: Initialize: Pseudo-label set  $\mathcal{P} \leftarrow \emptyset$ 
2: for each image  $I \in \mathcal{D}$  do
3:   Collect Predictions:  $\mathcal{A}(I) = \bigcup_{M_i \in \mathcal{M}} M_i(I)$ 
4:   Filter by Confidence: Retain detections with  $s(a) \geq \tau$ 
5:   Group by Class:  $\mathcal{A}^*(I, c) = \{a \mid \text{class}(a) = c\}$ 
6:   for each category  $c$  do
7:     Sort by Confidence (highest first)
8:     IoU Clustering: Initialize empty clusters
9:     for each detection  $a_i \in \mathcal{A}^*(I, c)$  do
10:      if  $a_i$  overlaps ( $\text{IoU} \geq \theta$ ) with an existing cluster then
11:        Assign  $a_i$  to the cluster
12:      else
13:        Create a new cluster with  $a_i$  as reference
14:      end if
15:    end for
16:    Soft Weighted Box Fusion (WBF):
17:    for each cluster  $\mathcal{C}_j$  do
18:      if detections from  $\geq m$  distinct models then
19:        Compute fused bounding box:

```

$$B^* = \sum_{a_i \in \mathcal{C}_j} w_i B_i, \quad w_i = \frac{s_i}{\sum s}$$

```

20:      Compute base confidence:  $s_{\text{base}} = \max(s_i)$ 
21:      Compute spread:

```

$$\text{spread} = 1 - \frac{1}{K} \sum_{j=1}^K \text{IoU}(B_j, B^*)$$

```

22:      Compute consensus factor:

```

$$\text{cf} = \exp(-\alpha \cdot \text{spread}) \cdot [1 + \beta \cdot (\text{num_models} - 2)]$$

```

23:      Compute soft confidence:  $s^* = s_{\text{base}} \cdot \text{cf}$ 
24:      if  $s^* \geq \tau_f$  then
25:        Add  $(B^*, s^*)$  to  $\mathcal{P}$ 
26:      end if
27:    end if
28:  end for
29: end for
30: end for
31: Clip Bounding Boxes to valid image dimensions
32: Output: Final soft pseudo-labels  $\mathcal{P}_{\text{Ensemble}}$  in COCO format

```

Qualitative Evaluation of Ensemble-Based Pseudo-Labels. Figure 9 visualizes the spatial alignment between our ensemble-fused pseudo-labels and manually verified ground-truth annotations. Solid cyan boxes denote fused predictions, while dashed magenta boxes represent ground truth. The close correspondence across diverse scenes, despite occlusions, deformation, and clutter, demonstrates the high localization quality and robustness of our pseudo-labeling pipeline, reinforcing its suitability for scalable semi-supervised training.



Figure 9: Qualitative comparison of fused pseudo-labels against ground-truth annotations on sample waste-sorting images. Solid cyan boxes depict our ensemble-fused detections, while dashed magenta boxes show manually-verified ground truth. Across four diverse scenes—varying object shapes, occlusions, and background clutter—the fused boxes closely align with the true annotations, demonstrating the robustness and high localization accuracy of our ensemble-based pseudo-labeling pipeline.

D Fine-Tuning Initializations: Checkpoints and Pre-Training Datasets

All detectors were initialized from official author-released checkpoints in the cited codebases. Table 10 reports, for each fine-tuned model, the backbone, codebase, config/recipe ID, checkpoint filename, and the pre-training dataset(s) documented for that checkpoint. For Grounding DINO, we used the language-aligned Swin-T and Swin-B checkpoints and froze the text encoder during fine-tuning.

Model	Backbone	Codebase	Config ID	Checkpoint Filename	Pre-Training Dataset(s)
YOLO11 (L)	–	Ultralytics	model=yolo111.pt, imsz=640	yolo111.pt	COCO
RT-DETR (L)	CSPResNet-50	Ultralytics	model=rtddetr-1.pt, imsz=640	rtddetr-1.pt	COCO
DINO	ResNet-50	Official repository	DINO_4scale.py	checkpoint0033_4scale.pth	COCO
DINO	Swin-L	Official repository	DINO_4scale_swin.py	checkpoint0029_4scale_swin.pth	Objects365, COCO
DETA	ResNet-50	Official repository	deta_ft	adet_2x_checkpoint0023.pth	COCO
DETA	Swin-L	Official repository	deta_swin_ft	adet_swin_ft.pth	Objects365, COCO
Co-DETR	ResNet-50	Official repository	co_dino_5scale_r50_lx.py	co_dino_5scale_r50_lx_coco.pth	COCO
Co-DETR	Swin-L	Official repository	co_dino_5scale_swin_large_l6_e_o365tococo.py	co_dino_5scale_swin_large_l6_e_o365tococo.pth	Objects365, COCO
Grounding DINO	Swin-T	MMDetection	grounding_dino_swin-t_finetune_l6xb2_lx_coco.py	grounding_dino_swin-t_finetune_l6xb2_lx_coco_20230921_152544-5f234b20.pth	Objects365, GoldG-VQA, Cap4M
Grounding DINO	Swin-B	MMDetection	grounding_dino_swin-b_finetune_l6xb2_lx_coco.py	grounding_dino_swin-b_finetune_l6xb2_lx_coco_20230921_152544-5f234b20.pth	COCO, Objects365, GoldG-VQA, Cap4M, Open Images, ODinW-35, RefCOCO

Table 10: Initialization details for fine-tuning on ZeroWaste-f. All initializations use official author-released checkpoints; the table reports the exact config/recipe and checkpoint used. Pre-training dataset(s) are as documented by the authors for each checkpoint (Grounding DINO uses language-aligned checkpoints with a frozen text encoder).